



Enhanced ethernet congestion management scheme for multicast traffic

Hela Mliki, Lamia Chaari, Lotfi Kamoun, Bernard Cousin

► To cite this version:

Hela Mliki, Lamia Chaari, Lotfi Kamoun, Bernard Cousin. Enhanced ethernet congestion management scheme for multicast traffic. Transactions on emerging telecommunications technologies, Wiley-Blackwell, 2016, 27, pp.1563 - 1579. 10.1002/ett.3097 . hal-01427027

HAL Id: hal-01427027

<https://hal.archives-ouvertes.fr/hal-01427027>

Submitted on 5 Jan 2017

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Enhanced Ethernet Congestion Management Scheme for Multicast Traffic

HELA MLIKI, LAMIA CHAARI, LOTFI KAMOUN, BERNARD COUSIN

LETI Laboratory, Engineering School of Sfax (ENIS), University of Sfax, 3038 Sfax, Tunisia

IRISA/University of Rennes 1, Campus de Beaulieu, 35042 Rennes cedex, France

mliki.hela@gmail.com, lamia.chaari@tunet.tn, lotfikamoun2@gmail.com , Bernard.Cousin@irisa.fr

Abstract

The Quantized Congestion Notification (QCN) is a Layer 2 congestion control scheme for Carrier Ethernet data center networks. The QCN has been standardized as an IEEE 802.1Qau Ethernet Congestion Notification standard. This paper reports the results of a QCN study with multicast traffic and proposes an enhancement to the QCN. In fact, in order to be able to scale up, the feedback implosion problem has to be solved. Therefore, we resorted to the representative technique, which uses a selected congestion point (i.e., the overloaded queue in a switch), to provide timely and accurate feedback on behalf of the congested switches in the path of multicast traffic. This paper evaluates the rate variation, the feedback overhead, the loss rate, the stability, the fairness, and the scalability performance of the standard QCN with multicast traffic and the enhanced QCN for multicast traffic. This paper also compares their performance criteria. The evaluation results show that the enhanced proposition of the QCN for multicast traffic gives better results than the standard QCN with multicast traffic. Indeed, the feedback implosion problem is settled by decreasing remarkably the feedback rate.

Index Terms – Ethernet Congestion Management scheme, Multicast traffic, IEEE 802.1Qau standard, QCN.

1 Introduction

THE IEEE 802.1 standards committee seeks to enable Data Center applications by using the Ethernet as an infrastructure [1], [2], [3]. To this end, the original set of Ethernet LAN technologies needs to support new capabilities to deliver enhanced services. This evolving set of Ethernet technologies is called the Carrier Ethernet [4], [5], [6].

In order to manage traffic overhead and improve QoS, IEEE 802.1Qau proposes Quantized Congestion Notification (QCN) as a scheme to manage congestion for a Carrier Ethernet network [7], [8]. The QCN relies on an end-to-end congestion notification to control congestion at the Layer 2 network. The switch that experiences a queue overload sends feedback notification frames over the network toward the source. The source has to adjust its data sending traffic according to the received feedback frames. The source is, then, called the reaction point (RP) and the switch that experiences congestion is called the congestion point (CP) [9], [10], [11].

Multicast is a communication mode that distributes one copy of the data from a source to an address group to be received by multiple destinations sharing this address. The group address is a multicast IP address in the IP network. However, it is mapped to a MAC group address in the Ethernet network [12].

When congestion occurs, the RP may face the feedback implosion problem, defined by a significant number of feedback frames, which may be returned by overloaded CPs for each copy of a multicast data frame sent.

Congestion control at Layer 2 using the QCN scheme has been studied extensively in the context of unicast traffic [13], [14], [15], [16], [17], [18], [19], [20]. However, the control of multicast traffic was not the object of the QCN standard [9]. In addition, there is a dearth of studies for the QCN in the case of multicast traffic. Nevertheless, providing a congestion control mechanism is critical in enabling multicast traffic for a Carrier Ethernet network. Therefore, we opted to focus on the standard QCN with multicast traffic in this paper. The proposed scheme builds a scalable congestion control multicast data link mechanism for a Carrier Ethernet without a feedback implosion problem. The feedback implosion problem is defined by an important number of feedback frames, which may be generated by overloaded CPs. Our enhancement proposition can avoid the feedback implosion problem to a great extent. One of our previous works compared QCN performance for multicast traffic with that for multiple unicast traffic [21]. While our previous work proposed to solve the feedback implosion problem by setting a high *Qeq* threshold value (for instance 75% of the queue capacity) [21], the proposed enhancement for QCN, in this paper, could solve the feedback implosion problem at low *Qeq* threshold values (25% of the queue capacity). Thus, our proposition for enhancement could fit the recommendation of the QCN standard (i.e., set the *Qeq* threshold at 25% of the queue capacity) [9], [22].

The following performance criteria for the QCN scheme for multicast traffic were addressed:

- i) Feedback overhead: how many feedback frames are generated when congestion is detected.
- ii) Loss rate: how many frames are dropped when a queue capacity is exceeded.
- iii) Stability: how the sending rate and the queue length fluctuate.
- iv) Fairness: how the QCN multicast traffic shares the bandwidth among sources.
- v) Scalability: how the QCN behaves when the number of switches and the multicast groups along a path increases. Indeed, a source needs to receive feedback frames from the CPs to determine the network traffic status in order to adjust its data rate accordingly. However, when the number of CPs increases, the multicast source can face a feedback implosion problem, which eventually results in a performance degradation.

The contribution of this paper is thus as follows: it first evaluates the QCN performance for multicast traffic in terms of feedback overhead, loss rate, stability, fairness and scalability through simulations. It also proposes an enhancement of the standard

QCN for multicast traffic using the representative technique. Then, it compares the QCN with multicast traffic to the enhanced QCN for multicast traffic.

The rest of the paper is organized as follows. The QCN congestion control scheme for Carrier Ethernet is presented in section 2. In section 3, we describe our proposition to enhance the standard QCN for multicast traffic in order to improve the network performance. In section 4, we describe the settings for our study, as well as the performance criteria used to evaluate the QCN with multicast and the enhanced QCN for multicast traffic. We report and discuss our findings in section 5. Finally, section 6 presents the conclusions drawn from our study.

2 Background

This section presents an overview of the QCN scheme according to [9].

The QCN monitors the queue utilization by requiring a queue length threshold (Qeq) at the output queues of the switch. When the queue length ($Qlen$) is beyond the threshold (Qeq), the queue manifests an indication that congestion is building up and a feedback frame is sent to the source to adjust its sending rate.

The Congestion Point (CP) and Reaction Point (RP) are the main parts of the QCN scheme to control congestion.

- i) Congestion Point (CP): It detects congestion by monitoring the switch queue length. The aim is to prevent the queue length from exceeding the queue threshold Qeq . The CP signals congestion by generating a feedback frame addressed to the source of the sampled frame that causes congestion. Indeed, once a data frame is sampled, the CP measures the congestion level on the link. Therefore, the CP computes a congestion measure Fb . This measure will be held into a feedback frame to notify the source about congestion.
- ii) Reaction point (RP): It is associated with a source to adjust the sending flow rate. The rate is decreased when a feedback frame is received. However, the rate is increased when the RP deduces that there is available bandwidth.

When the computed Fb value at the CP is positive, no congestion is detected by the switch and no feedback frames are sent to the RP. The RP, then, infers that it could increase its transmission rate. When the Fb value is negative at the CP, a feedback frame is generated, then, the RP decreases its rate.

2.1 Congestion Point

The RP defines how the source rate is adjusted, while the CP defines how the congestion measure Fb is updated. For every received data frame, the CP checks the queue occupancy. The CP detects congestion when the computed Fb value is negative. Then, the CP notifies the source a congestion status by generating a feedback frame. A feedback frame is sent to the sampling frame that caused congestion so that the system converges to fairness. The feedback frame carries the Fb value, which is used to communicate the switch queue state to the RP. The Fb value is quantized to 6 bits. Thus, the maximum quantized Fb value (Fb_{max}) is equal to 63.

The Fb is updated as follows:

$$Fb = -(Qoff + w \times Qdelta) \quad (1)$$

Here, w is a non negative constant, chosen to be 2 in the standard QCN [9].

$Qoff$ represents the queue size excess while $Qdelta$ represents the rate excess; they are defined as follows:

$$Qoff = Qlen - Qeq \quad (2)$$

$$Qdelta = Qlen - Qold \quad (3)$$

Here, $Qlen$ denotes the instantaneous queue size. However, $Qold$ denotes the queue size when the last feedback message was generated.

The main objective of the QCN is to prevent, as much as possible, the queue from building up to the point at which a frame has to be dropped, therefore, it uses a threshold (Qeq).

2.2 Reaction Point

The QCN adapts the source rate to the existing network status; it increases its sending rate if the network appears to be free of congestion (i.e., when no feedback frame is received) and decreases the source rate if the network suffers from congestion (i.e., when a feedback frame is received).

Let CR denote the current sending rate of the source data traffic, and TR denote the RP sending rate just before receiving a feedback frame. The RP aims to keep the CR data transmission rate from the source below the TR . The RP decreases its data traffic rate when a feedback frame is received.

When the RP receives a feedback frame, it deduces that congestion occurs. Therefore, it decreases its current rate (CR) and updates its target rate (TR) as follows:

$$CR = CR \times (1 - Gd \times |Fb|) \quad (4)$$

$$TR = CR \quad (5)$$

Here, the constant Gd is selected so that $Gd \times |Fb_{max}| = 1/2$. Then, the current rate can decrease by 50% [9].

The multiplicative decrease is expected to reduce an overload at the queue in the CP, the RP is expected to be able to increase its rate afterwards. This helps to recover some of the lost bandwidth. When no feedback frame is received, the QCN performs the following increase phases: Fast Recovery (FR), Active Increase (AI) and Hyper Active Increase (HAI).

2.2.1 Fast recovery

When no feedback frame is received, the RP performs five cycles of FR to increase the sending rate. To compute the duration of each of the five cycles of the FR phase, the RP uses a Byte Counter, which counts the number of bytes transmitted and a Timer, which times the rate increase. During each FR cycle, the RP transmits 150 *Kbytes*. The timer completes one cycle with T duration ($T = 10\text{ ms}$)

At the end of each cycle, the TR does not change, while the CR is updated as follows:

$$CR = \frac{1}{2} \times (TR + CR) \quad (6)$$

2.2.2 Active Increase

After performing the five cycles of the FR phase and no feedback frame is received, the RP deduces that there is an available bandwidth. It switches to perform the next AI phase where it increases the Current Rate (CR) more than the previous phase. The AI phase also uses a Byte Counter and a Timer each of which is equal to one cycle. The RP transmits 75 *KBytes*. The timer completes $T \div 2$ duration ($T = 5\text{ ms}$). Then, at the end of the cycle, the RP updates the TR and CR as follows:

$$TR = TR + R_{AI} \quad (7)$$

$$CR = \frac{1}{2} \times (CR + TR) \quad (8)$$

Here, R_{AI} is a constant selected to be 5 *Mbps* in the QCN standard [9].

2.2.3 Hyper Active Increase

At the end of the AI phase and if no feedback frame is received, the RP deduces that there is available bandwidth. It switches to perform the next HAI phase where it increases the Current Rate (CR) substantially. The RP increases the TR and CR as follows:

$$TR = TR + i \times R_{HAI} \quad (9)$$

$$CR = \frac{1}{2} \times (CR + TR) \quad (10)$$

Here, i is the number of HAI cycles, selected to be equal to one and R_{HAI} is set to 50 *Mbps* in the QCN standard [9].

When a feedback frame is received during an increase phase (FR or AI or HAI), the increase phase is cancelled: the Byte Counter and the Timer are set to zero. Then, a multiplicative decrease is performed as it is described above.

3 The Enhanced QCN for Multicast Traffic

In this section, we present the key idea of the enhanced QCN for multicast traffic. The added operations at a congestion point and at the reaction point are detailed.

Our work is inspired from the representative technique. This technique was an early single rate multicast congestion control scheme defined in DeLucia et al.'s work [23]. PGMCC [24], TFMCC [25] and MDP-CC [26] are examples of well known schemes that use the representative technique. This technique defines a small set of multicast group members that can represent the congested multicast subtree. These group representative provide an immediate feedback packet, which can suppress any feedback from other group members, thus, preventing feedback implosion at the source. If a receiver never experiences congestion, or has its packet losses covered by a representative, it will never generate any feedback messages. These schemes were adapted to be implemented with TCP. Indeed, they all make use of the TCP throughput formula [27], which provides the receiver with the lowest estimate TCP throughput.

Our work leveraged the representative technique to boost the QCN scheme for multicast traffic without any major alteration of the QCN specification. The enhanced QCN for multicast traffic scheme proposes to define a selected congestion point among all the potential existing congestion points used to represent the congested multicast transmitted path to the destination. The selected congestion point provides immediate feedback frames, which can suppress feedback frames from other potential congestion points, thus, preventing feedback implosion at the source.

The scheme reacts to any new congestion a timely way by selecting a new representative and discarding those that are no longer contributing to the feedback efforts.

The representative in our scheme is the CP that has the greatest feedback $|Fb|$ value. If a CP is selected as a representative, only the feedback frames from that CP are allowed to be generated. The other CPs will cancel the action of feedback frame generation. Then, The RP compares between its Fb value (computed from previous feedback frames) and the Fb value that it receives in a feedback frame from a CP. The Fb at the RP is updated when its value is lower than the Fb value received in the feedback frame. This makes the RP aware about the largest congestion in the multicast tree of its transmitted flow in order to decrease its transmission rate. In addition, QCN has increasing rate phases (i.e., AI, FR, and HAI) used to recover the data rate that could have been lost during the last rate decrease episode and to grab extra available bandwidth. The aim of our enhancement is to receive feedback from one representative CP on each path of the multicast session of the source traffic in order to avoid the feedback implosion problem. This representative CP is selected to be the one that suffers from the worst congestion case.

Moreover, when congestion is detected the congestion point sends a feedback frame. The feedback frame carries the Fb value, which is used to communicate to the RP the switch queue state. The Fb value is quantized to 6 bits. Thus, the maximum quantized Fb value (Fb_{max}) is equal to 63 (this is defined by the standard). Consequently, the maximum congestion measure that a reaction point could receive is 63. Therefore, when the feedback value at the reaction point reaches its maximum value (i.e., 63), it should be reset to zero in order to reset the selection process of the rep-

representative congestion point. Indeed, when the feedback value at the reaction point reaches the maximum value, there is no worst congestion measure of feedback to be quantized at the CP and sent back to the RP. Thus, the reaction point that received the maximum quantized feedback measure could no longer trigger the congestion point (i.e., the representative CP) to send a feedback frame: in such a case, the computed quantized Fb value at the congestion point could be changed to a less value than the one received in the data frames from the reaction point. Therefore, in order to select a new representative CP and trigger a new feedback frames, the Fb is reset to 0 when it reaches 63 at the reaction point.

3.1 The Proposed Enhancement at the Reaction Point

Operations at the reaction point as defined by the standard [9] remain unchanged. However, some other operations have been added to carry out the enhancement. These operations make the reaction point responsible for distributing the current representative set.

Each time a reaction point receives a feedback frame, it compares the new feedback value $|Fb|$ with the previously received one. If the new received feedback value is greater than the previous one, the data frame source will hold the new value of the feedback. Otherwise, the reaction point keeps the feedback field of the transmitted data frame unchanged.

It is obvious that our proposition requires a field in the source data frame to hold the feedback $|Fb|$ value of the representative. Moreover, when the feedback value at the source reaches its maximum value (i.e., 63), it should be reset to zero in order to reset the selection representative process. Indeed, when the feedback value at the RP reaches the maximum value, there is no worst congestion measure of feedback to be quantized.

Algorithm 1 describes the additional operations performed when a feedback frame is received.

3.2 The Proposed Enhancement at the switch

The operations at a switch as defined by the standard [9] remain unchanged. However, some other operations have been added to achieve the enhancement. These operations consist in defining which congestion point is the representative.

When the congestion point detects congestion, it sends a feedback frame only if its computed feedback $|Fb|$ value is greater than the feedback value held in the received data frame from the source.

Thus, a congestion point designates itself as a representative when it has the greatest feedback $|Fb|$ value. Indeed, the feedback Fb value defines the measure of congestion, it captures a combination of queue size excess ($Qoff$) and rate excess ($Qdelta$). When the Fb is negative, it means that the queue is overloaded and a feedback frame is generated to notify the source about the state of congestion.

When the CP has the greatest feedback $|Fb|$ value, it can infer that it has the most overloaded queue among all queues in the path of the transmitted multicast traffic.

Algorithm 1 Pseudo-code of the proposed enhancement at the reaction point

Variables:

Fb the computed feedback at the received feedback frame from a CP to the RP;

\hat{Fb} the saved feedback value of the last received feedback frame from a CP to the RP;

Initialization:

$\hat{Fb} \leftarrow 0$;

When a feedback frame is received:

Get the Fb from the feedback frame;

if ($\hat{Fb} < Fb$) **then**

$\hat{Fb} \leftarrow Fb$;

endif

Adjust the data rate (\hat{Fb});

if ($\hat{Fb} == 63$) **then**

$\hat{Fb} = 0$;

endif

Therefore, it is the only one who can generate feedback frames. This is achieved by comparing the congestion measure $|Fb|$ between queues of potential congestion points.

Algorithm 2 describes the additional operations performed when a data frame is received.

4 Evaluation and simulation of enhanced QCN for multicast traffic

This section reports the performance of our proposition to enhance the QCN for multicast traffic through simulations and measurements. In addition, a comparison between the QCN with multicast traffic and the enhanced QCN for multicast traffic is achieved.

In this evaluation the following metrics are our major concern:

- i) Feedback overhead: computes the rate of feedback generated by the CP.
- ii) Loss rate: computes the rate of dropped frames by the CP.
- iii) Stability: computes the standard deviation (StdDev) of the rate CR at the RP and the queue length deviation at the CP.
- iv) Fairness: computes the fairness index [28] to measure fairness among the sources. The results of the fairness index are always a number between 0 and 1, with 1 representing the greatest fairness.

Algorithm 2 Pseudo-code of the proposed enhancement at the switch

Variables:

Fb the computed feedback at a CP;

\hat{Fb} the feedback value holds in a received data frame;

representative indicates whether the CP is the representative CP for multicast traffic;

Initialization:

representative \leftarrow false;

When a data frame is received:

Get the \hat{Fb} from the data frame;

Compute the Fb of the CP; /* according to Eq.1 */

if ($|Fb| \geq \hat{Fb}$) **then**

representative \leftarrow **true**;

else

representative \leftarrow **false**;

endif

if ($(Fb < 0) \ \&\& \ (\textit{representative} \text{ is } \textbf{true}))$ **then**

send feedback frame ($|Fb|$);

endif

- v) Scalability: computes the feedback rate, the loss rate, the stability and the fairness performance criterion when the number of CPs and multicast groups increases.

Carrier Ethernet network architecture could have many issues like congestion control, bridge loop, admission control, energy saving, etc. Since we address the congestion control issue in a Carrier Ethernet network, we selected topologies and scenarios that could exist in Carrier Ethernet network architecture and implement congestion. The chosen topologies are broadly used in the literature to study the network congestion state. In our study we want to show how the QCN congestion control scheme behaves. Indeed, for each issue there are topologies that are well known and well adapted to them. Many topologies and congestion scenarios are proposed in the literature; however, these scenarios should lead to the same analysis. Therefore, we selected only a small number of congestion scenarios, which are simple to understand and easy to analyse.

Thus, two topologies are used for the performance evaluation: a star topology defined as in Figure 1 and a multi-link topology defined as in Figure 13. The topology described in Figure 1 is used to study the feedback overhead, loss rate, stability and fairness, while the scalability is studied within the topology described in Figure 13.

A star topology defined as in Figure 1 has a single switch. However, this single switch holds two potential congested output queues: the queue that stores traffic for receiver R1 and the queue that stores traffic for receiver R2. Our enhanced proposition for the QCN has to select the appropriate congested queue to make the adequate adjustment at the RP and to not flood the source with feedback frames. The key idea is to select the queue that computes the greatest $|Fb|$ value. Since the $|Fb|$ defines a measure of congestion, the queue that has the greatest feedback $|Fb|$ value can infer that it is the most overloaded queue among the all queues in that congestion point. Thus, it is selected as a representative as defined in section 3. It is worth knowing that the standard QCN defines a Congestion Point Identifier (CPID), which is hold in the feedback frame. CPID field must be unique across the network and it is used to identify a congestion entity (i.e., a queue) [9]. Therefore, a selected representative is identified by its CPID. Thus, if there are two identical $|Fb|$ values from two different queues and are equal to the Fb value in the data frame, only the queue with the CPID that matches with the selected representative will generate a feedback frame.

We used the OMNeT++ simulator for this performance evaluation. There is one queue per switch output port. We used drop tail queues with FIFO scheduling. All queues have the same size and their total size is equal to *100 frames*; each is *1500 bytes* long. Our network used Ethernet links with a capacity equal to *1 Gbit/s*. The initial value of the *CR* and *TR* are set to the transmission rate of the Ethernet interface (*1 Gbit/s*). There are six sources and each one sends traffic at *200 Mbit/s* with a constant UDP frame size and a constant UDP frame generation time.

We used UDP-based traffic as a yardstick for our case of study. Multicast traffic is generally based on UDP. The UDP uses no congestion control mechanisms. Therefore, congestion control at Layer 2 is a valuable alternative in such a case.

We have studied a case where the QCN is used as a congestion control scheme in a network with multicast traffic, and a case where the QCN is enhanced by using

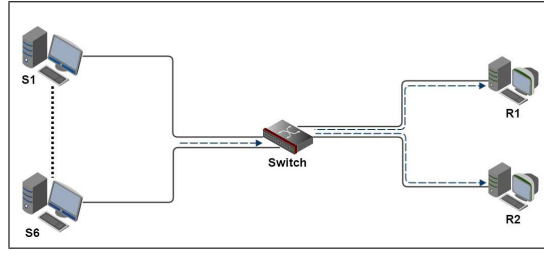


Figure 1: Star Topology

the representative technique to better handle multicast traffic in a Carrier Ethernet network.

Although we started with the objective to solve the QCN feedback implosion problem, our proposition of the enhancement of the QCN for multicast traffic revealed several additional advantages. It has reduced the loss rate and improved the scalability.

4.1 Rate variation

Figures 2 and 3 plot the rate variation during simulation time of each of the six multicast flows respectively when $Qeq=25 \text{ frames}$ and $Qeq=75 \text{ frames}$ in the case of the QCN with multicast traffic. Figures 4 and 5 plot the rate variation of each of the six multicast flows during simulation time when $Qeq=25 \text{ frames}$ and $Qeq=75 \text{ frames}$, respectively, in the case of the enhanced QCN for multicast traffic.

As there is no congestion at the beginning of the simulation, the RP sends traffic at a maximum CR. Then, as it receives feedback frames, the CR is decreased.

The CR_{mean} represents the black line where the CR of each flow should converge ($CR_{mean} = 1 \text{ Gbit/s} \div 6$) during simulation. In both schemes of the QCN (i.e., the QCN and the enhanced QCN), the CR fluctuates over the CR_{mean} when the Qeq is increased (Figures 3 and 5).

Figure 6 shows the RP transmission rate average for both multicast cases of the standard QCN and the enhanced QCN. Although the main objective is to reduce the feedback rate in order to avoid the feedback implosion problem, the enhanced scheme for the QCN can also improve the CR transmission rate compared to the standard QCN when multicast traffic is considered. Thus, thanks to the representative selection technique used with the enhanced QCN, which could eliminate some unnecessary feedback frames, the source could adjust its transmission rate (CR) less severely compared to the standard QCN. Indeed, with the enhanced QCN, when $Qeq=25 \text{ frames}$ the CR is increased by 12.48% of the QCN with multicast result.

4.2 Feedback Overhead

Figure 7 shows the feedback rate according to different values of the Qeq threshold in the case of the QCN with multicast and the case of the enhanced QCN for multicast traffic. We note, in both cases, that when the Qeq threshold value decreases, the generation of feedback frames increases. Indeed, low Qeq threshold values can easily be exceeded and then a congestion notification occurs. In addition, the enhanced QCN

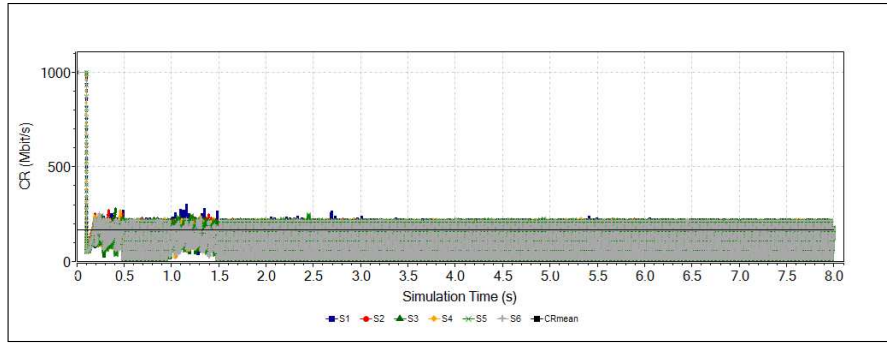


Figure 2: CR when Qeq=25 frames in the case of the QCN with multicast traffic

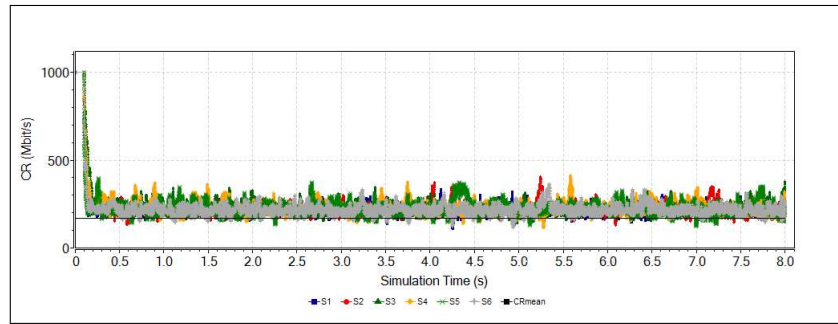


Figure 3: CR when Qeq=75 frames in the case of the QCN with multicast traffic

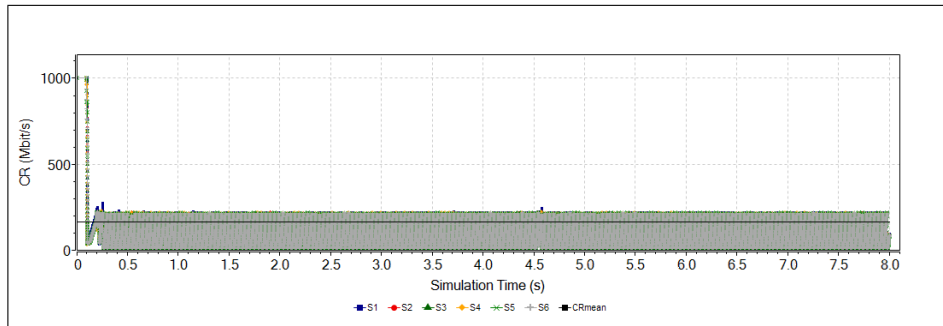


Figure 4: CR when Qeq=25 frames in the case of the Enhanced QCN for multicast traffic

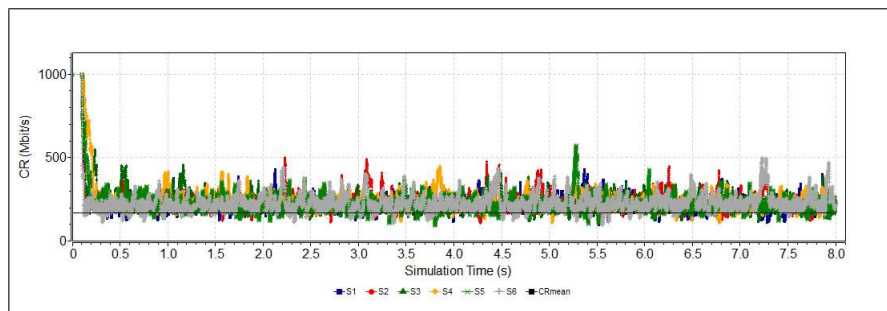


Figure 5: CR when Qeq=75 frames in the case of the Enhanced QCN for multicast traffic

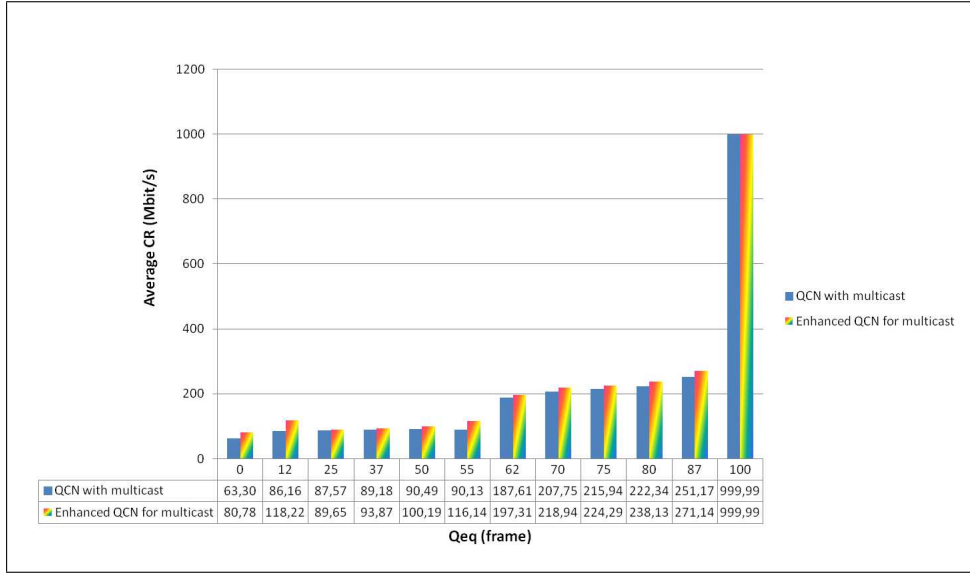


Figure 6: Average CR with a star topology

for multicast traffic proposed solution has succeeded in reducing the feedback rate. For example, in the case of the QCN with multicast traffic, when the $Qeq=75$ frames the feedback rate is equal to 3.8% of the total frame stream; when the $Qeq=50$ frames this rate is equal to 11.98% of the total frame stream; and when the $Qeq=25$ frames it is equal to 13.16%. However, in the case of the enhanced QCN for multicast traffic, when the $Qeq=75$ frames the feedback rate is equal to 2.27% of the total frame stream; when the $Qeq=50$ frames this rate is equal to 5.63%; and when the $Qeq=25$ frames it is equal to 8.04%. This means that we succeeded in making a reduction of -40.26% of the QCN with multicast result when the $Qeq=75$ frames, -53% when the $Qeq=50$ frames, and -38.9% when the $Qeq=25$ frames. It is then obvious that the enhanced QCN for multicast traffic can decrease the feedback overhead significantly.

4.3 Loss Rate

Figure 8 shows the frame loss rate according to different Qeq threshold values in the case of the QCN with multicast and the case of the enhanced QCN for multicast traffic.

We note that as the Qeq threshold value increases, the loss rate also increases in both the QCN with multicast and the enhanced QCN for multicast traffic. This is because the low Qeq threshold value leaves a safety margin for burst arrivals of new flows. That is why it has a lower drop rate than those of a high Qeq threshold. For example, in the case of QCN with multicast traffic, when the $Qeq=25$ frames the loss rate is equal to 0% of the total frame stream; when the $Qeq=50$ frames this rate is equal to 3.15% of the total frame stream; and when the $Qeq=75$ frames it is equal to 26.59% of the total frame stream. However, in the case of the enhanced QCN for multicast traffic, when the $Qeq=25$ frames the loss rate is equal to 0% of the total frame stream; when the $Qeq=50$ frames this rate is equal to 2.17%; and when the $Qeq=75$ frames it is equal to 29.98%. On the one hand, when the Qeq threshold values are low, the enhanced QCN for multicast traffic has better results in terms of loss rate than the case of the

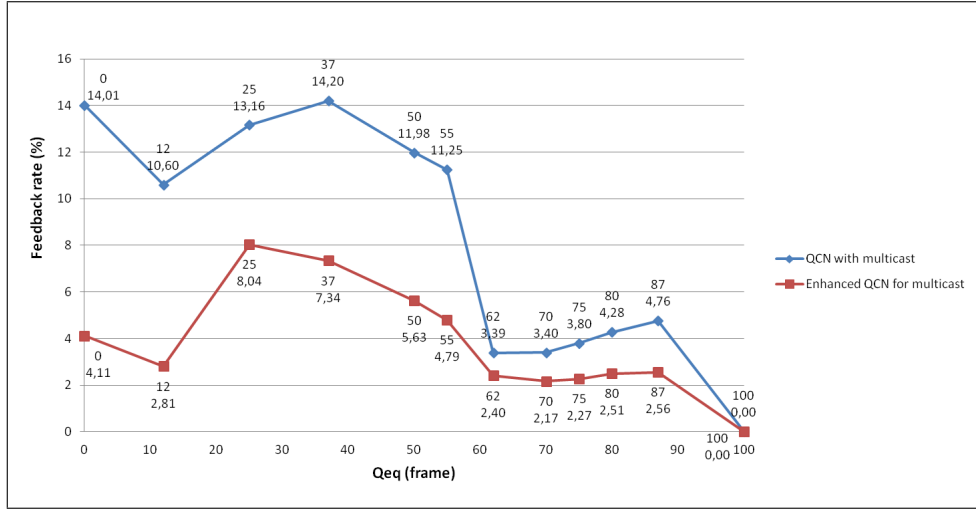


Figure 7: Feedback rate with a star topology

QCN with multicast traffic. For example when the $Qeq=25$ frames the loss rate in both cases is equal to 0% of the total frame stream, whereas when the $Qeq=50$ frames the enhanced QCN for multicast succeeded in making a reduction of -31.11% of the QCN with multicast result. On the other hand, when the Qeq threshold values get higher, the enhanced QCN for multicast traffic has worse results in terms of loss rate than the QCN with multicast traffic. For example, when the $Qeq=75$ frames the enhanced QCN for multicast achieves an increase of 12.74% of the QCN with the multicast result. This could be explained by the lack of feedback frames in the enhanced QCN for multicast traffic case compared to the QCN with multicast traffic case. The lack of feedback frames prevents the adjustment of the transmission rate properly at sources and, then, increases the loss rate.

The standard recommendation for QCN is to set the Qeq threshold at low value because it allows a tight congestion control by sending feedback frames early. This enables an early adjustment of the transmission rate, provides a safety margin for burst arrivals of new flows and could decrease queue delay. Thus, this has motivated us to study an enhancement for the QCN in order to improve performance at low Qeq threshold values for multicast traffic. However, our performance results included also high Qeq threshold values to study the impact of our scheme at these high queue levels. Our scheme could occasionally have some better results for high Qeq threshold value (although we think that this is not a significant result), but our study objective is to have better results for low Qeq threshold values in order to fit the standard recommendation of the QCN.

4.4 Stability

The stability performance criteria characterizes the fluctuation magnitude of the system variables.

We study the stability of the sending rate CR and the queue length for both the QCN with multicast and the enhanced QCN for multicast traffic.

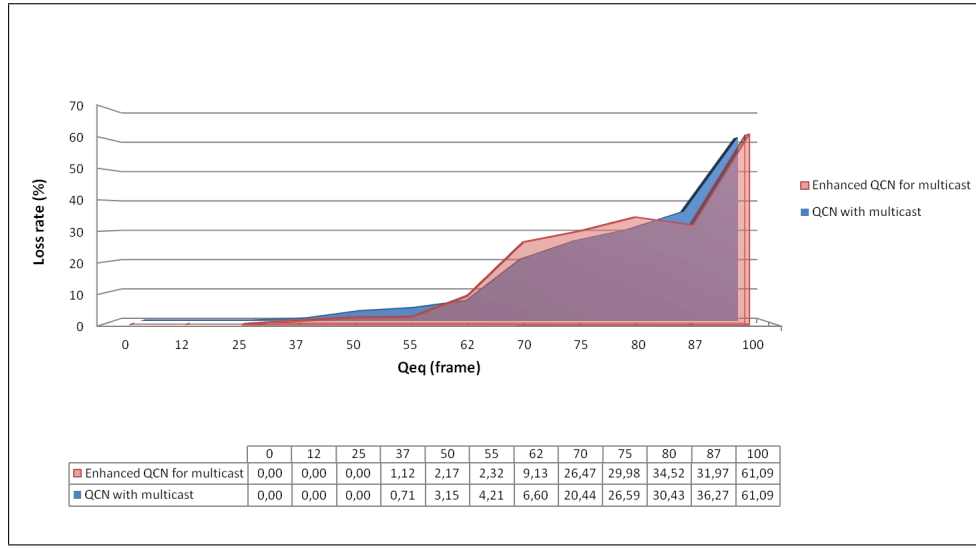


Figure 8: Loss rate with a star topology

Figure 9 shows the standard deviation (stdDev) of the CR at the RP for different Qeq threshold values in the both QCN schemes. We note that the CR experiences more fluctuations in the case of the enhanced QCN for multicast traffic than the case of the QCN with multicast traffic when the Qeq threshold value increases. However, the proposed solution of the enhanced QCN for multicast traffic succeeds in reducing this fluctuation when the Qeq threshold values are low. For example, in the case of the QCN with multicast traffic, when the $Qeq=75$ frames the stdDev of the CR is equal to 43.61 Mbit/s; when the $Qeq=50$ frames this stdDev is equal to 69.48 Mbit/s; and when the $Qeq=25$ frames it is equal to 67.65 Mbit/s. However, in the case of the enhanced QCN for multicast traffic, when the $Qeq=75$ frames the stdDev of the CR is equal to 58.56 Mbit/s; when the $Qeq=50$ frames this stdDev is equal to 65.19 Mbit/s; and when the $Qeq=25$ frames it is equal to 65.32 Mbit/s.

Figure 10 shows the mean queue length for different values of Qeq threshold for the both QCN schemes. Figure 11 plots the deviation of the mean queue length from the Qeq threshold in the case of the QCN with multicast traffic and the case of the enhanced QCN for multicast traffic.

The purpose is not to exceed the Qeq threshold while transmitting the source frames to their destinations. It is important to highlight that the mean queue length is under the Qeq threshold when the deviation value is negative, but it exceeds the Qeq threshold when the deviation is positive. If the mean queue length goes over the Qeq threshold it indicates a poor control of congestion. We find that, in the case of the enhanced QCN for multicast traffic, congestion is inadequately controlled in the case of high Qeq threshold compared to the QCN with multicast traffic. However, the proposed solution to enhance the QCN for multicast traffic could control congestion properly in the case of low Qeq threshold values.

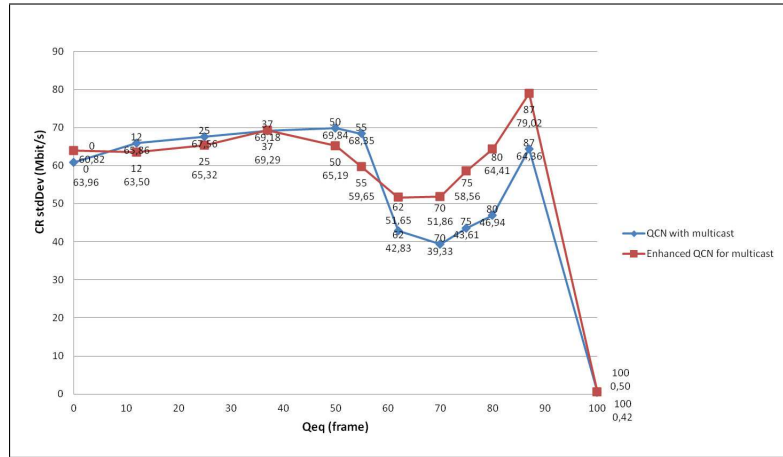


Figure 9: Standard deviation of the CR with different Qeq threshold values with a star topology

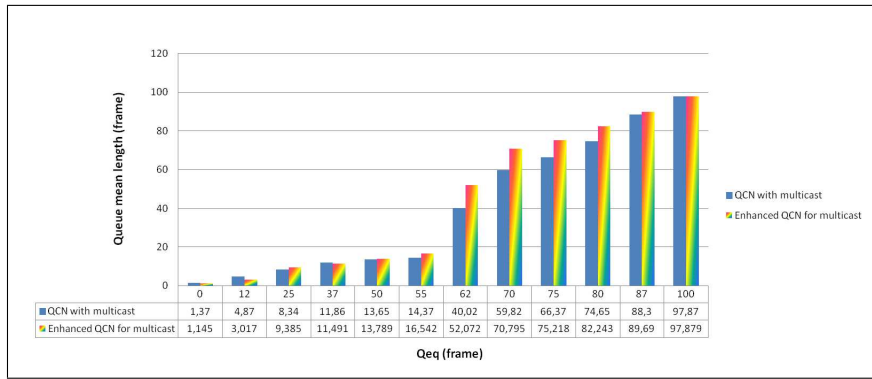


Figure 10: Queue mean length variation for different Qeq threshold values with a star topology

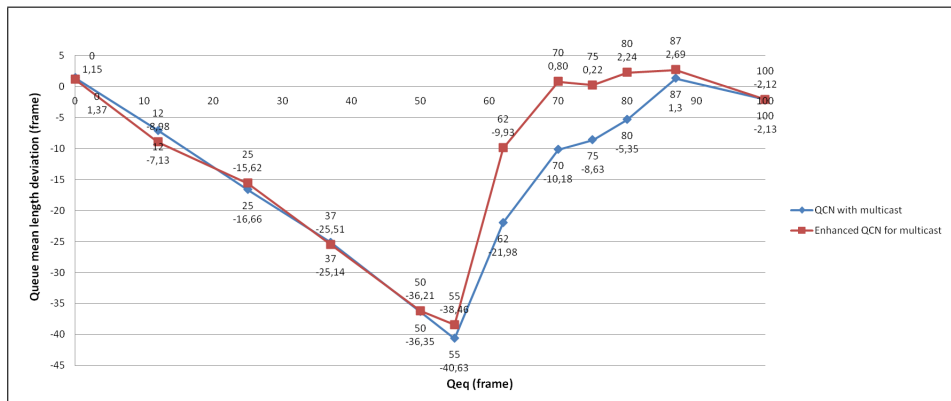


Figure 11: Deviation of the mean queue length from the Qeq thresholds with a star topology

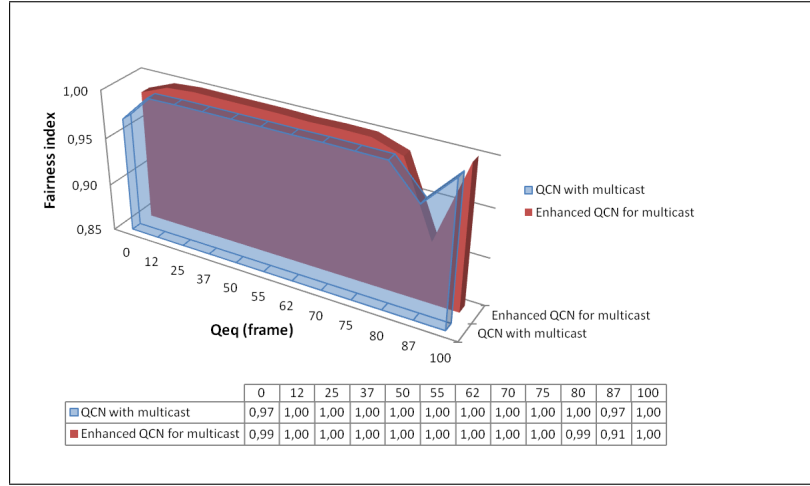


Figure 12: Fairness index with a star topology

4.5 Fairness

According to max-min fairness [29], network resources are allocated in such a way that the bit rate of a flow cannot be increased without decreasing the bit rate of a flow with a smaller bit rate.

Figure 12 plots the fairness index comparison between the QCN with multicast traffic and the enhanced QCN for multicast traffic for different Qeq thresholds. We find that the fairness indices are similar for both QCN schemes for most of the Qeq threshold values.

4.6 Scalability

The scalability performance helps to check the performance of the standard and the enhanced QCN when the network scales up. The QCN scheme could be implemented within a data center network. In a data center network we could get different congested link. We scale down the data center architecture into topology in Figure 13 to address congestion issue with multiple bottlenecks. Indeed, with this topology three potential bottleneck are involved. It goes without saying that increasing the number of bottleneck leads to same analysis.

Figure 13 shows a scenario with a multiple path used to study the scalability performance. This scenario implements a network with an increased number of CPs and multicast groups. With this scenario, it is expected to have heavy congestion that leads to feedback implosion.

Multicast transmission mode consists in sending a single copy of data traffic to a selected group of destination. We proposed a scenario with three multicast group address. Three receiving multicast groups (G1, G2, G3) are specified to the same number of sources at the reception step. In the first group (G1) the multicast flow goes through three potential CPs (Switch1, Switch2, Switch3), from six sources (S1, S2, S3, S4, S5, S6) to receivers (R5, R6). However, in the second group (G2) the multicast flow goes through two potential CPs (Switch2, Switch3), from six sources to receivers (R1,

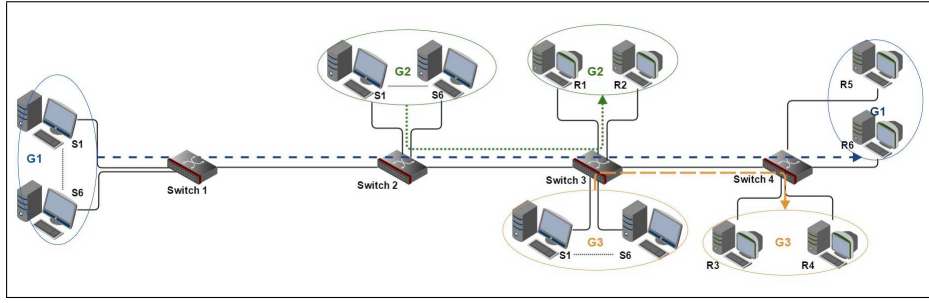


Figure 13: Multi-link topology with multiple potential bottlenecks

R2). Ultimately, in the third group (G3) the multicast flow goes through one potential CP (Switch3), from six sources to receivers (R3, R4). The QCN with multicast traffic and the enhanced QCN for multicast traffic are then studied through this scenario to compare their performance criteria. In QCN standard, congestion points will send a feedback notification to the RP. However, with our proposition of QCN enhancement, only one selected congestion point on each multicast path is allowed to send its feedback frame to the source. The challenge is to select the appropriate congestion point as a good representative of congestion level in order to not flood the source with feedback frames. Our enhanced QCN for multicast traffic scheme proposes to select one congestion point among all the existing congestion points. This selected congestion point represents the congestion level for its path on the multicast tree. The representative in our scheme is the CP that has the greatest feedback $|Fb|$ value. If a CP is selected as a representative, only the feedback frames from that CP are allowed to be generated. The other CPs cancel their feedback frame generation. The selected congestion point provides immediate feedback frames, which can suppress feedback frames from other potential congestion points, thus, preventing feedback implosion at the source.

4.6.1 Feedback Overhead in a Multi-Link Topology

Figure 14 illustrates the rate of the feedback notifications received at the RPs with different Qeq threshold values. This figure compares the feedback rate between the two QCN schemes. For example, in the QCN with multicast traffic case, when the $Qeq=25$ frames the feedback rate is equal to 87.92% of the total frame stream; when the $Qeq=50$ frames this rate is equal to 89.02% of the total frame stream; and when the $Qeq=75$ frames it is equal to only 4.24% of the total frame stream. However, in the enhanced QCN for multicast traffic, when the $Qeq=25$ frames the feedback rate is reduced by -41% of the QCN with multicast result; when the $Qeq=50$ frames this rate is reduced by -41.23% ; and when the $Qeq=75$ frames it is reduced by -60.14% . It is obvious that the proposed solution to enhance QCN for multicast traffic generates a smaller feedback rate compared to the QCN with multicast traffic.

4.6.2 Loss Rate in a Multi-Link Topology

Figure 15 compares the frames loss rate of the two QCN schemes for different Qeq threshold values. It can be deduced that the frame loss rate is high when the Qeq

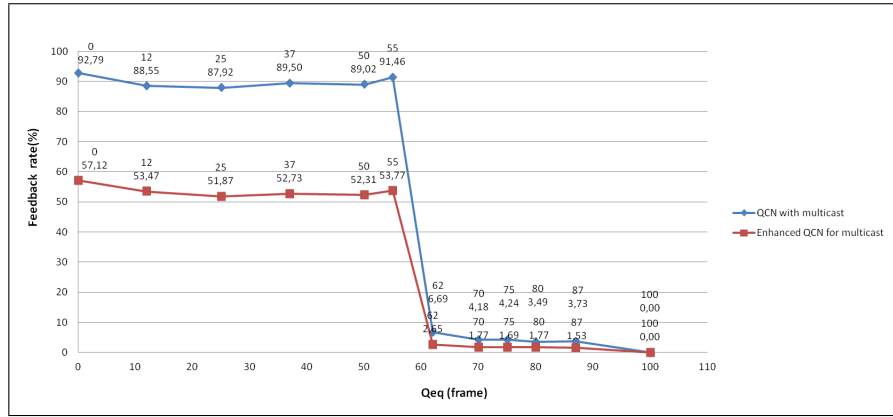


Figure 14: Feedback rate with a multi-link topology

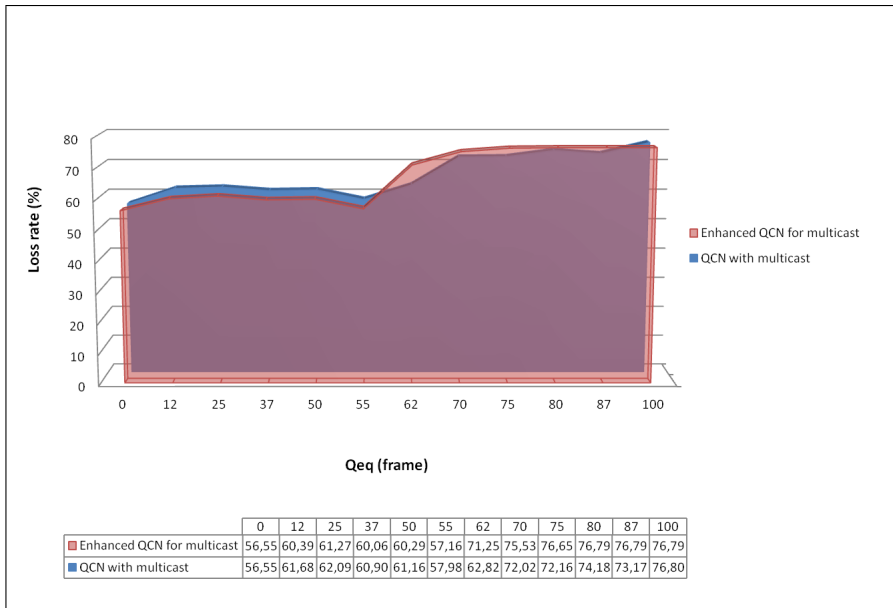


Figure 15: Loss rate with a multi-link topology

threshold values go up.

The enhanced QCN for multicast traffic decreased the loss rate compared to the QCN with multicast traffic in the case of low values of Qeq threshold. However, the loss rate of the enhanced QCN for multicast traffic gets increased compared to the QCN with multicast traffic when the Qeq threshold values go up.

The enhanced QCN for multicast traffic shows a low loss rate control, when the Qeq threshold values go up, compared to the QCN with multicast traffic. Indeed, there is already a low feedback rate with the QCN with multicast when the Qeq threshold value increases; with the representative selective technique used by the enhanced QCN for multicast traffic, the sources receive a poor notification information of congestion represented by a low number of feedback frames. Consequently, this does not allow the RP to adjust the sending rate adequately and thus a frame loss occurs when the Qeq threshold value increases.

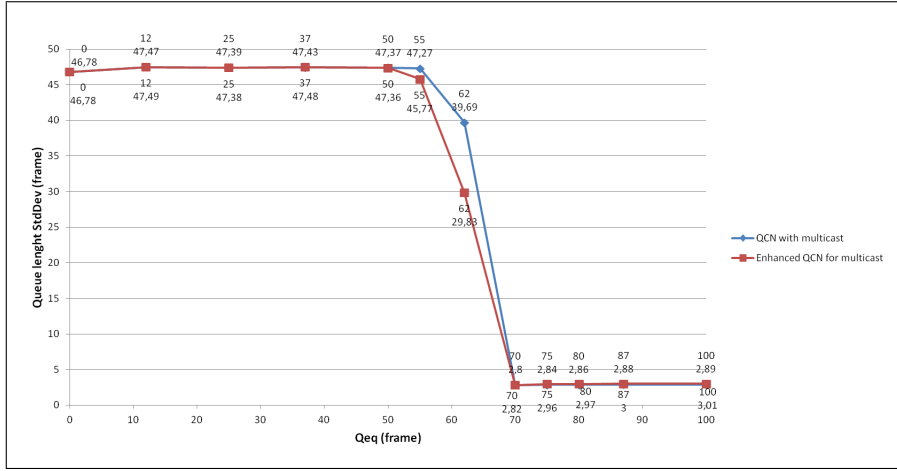


Figure 16: Queue length StdDev for different Qeq thresholds in switch 1 with the multi-link topology

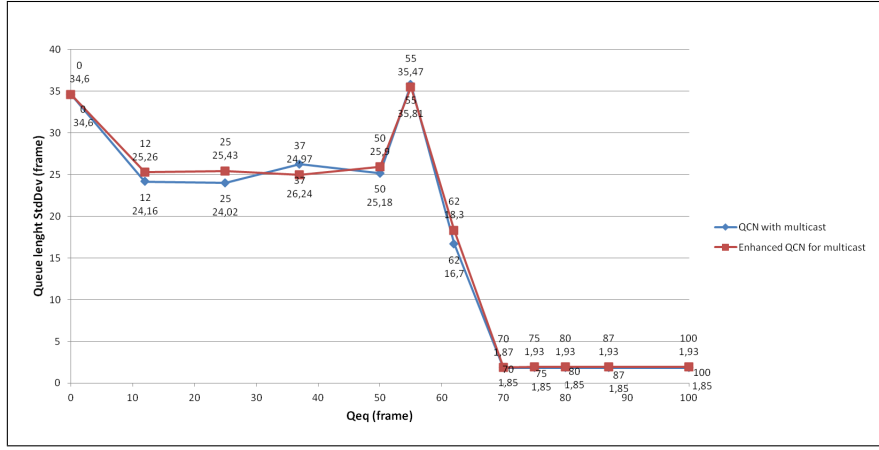


Figure 17: Queue length StdDev for different Qeq thresholds in switch 2 with the multi-link topology

4.6.3 Stability in a Multi-Link Topology

Figure 16 plots the queue length stdDev for different values of Qeq threshold at switch 1 of the multi-link topology. Figure 17 shows the queue length stdDev for different values of Qeq threshold at switch 2 of the multi-link topology. Figure 18, however, displays the queue length stdDev for different values of Qeq threshold at switch 3 of the multi-link topology.

We note that the StdDev of the queue length at CPs (Switch 1, Switch 2 and Switch 3) of the QCN with multicast traffic and the enhanced QCN for multicast traffic is almost similar for different Qeq threshold values.

4.6.4 Fairness in a Multi-Link Topology

Figure 19 compares the fairness index of the QCN with multicast traffic to the enhanced QCN for multicast traffic at G1 sources. Figure 20 displays the fairness index of the QCN with multicast traffic and the enhanced QCN for multicast traffic at G2 sources. Figure 21 shows the fairness index of the two QCN schemes at G3 sources.

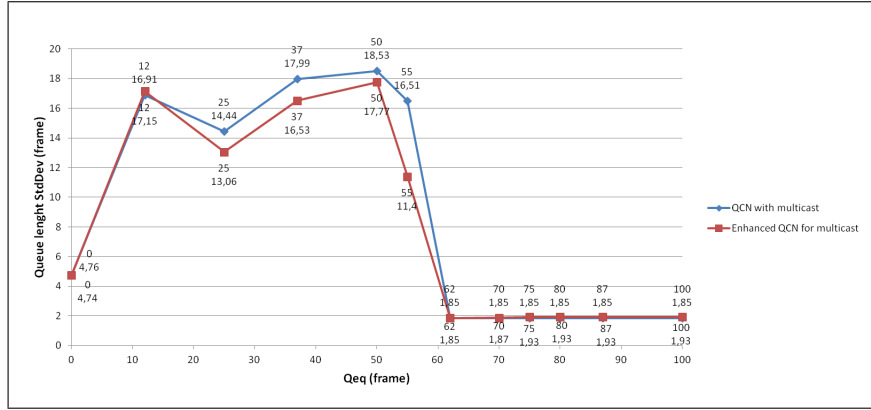


Figure 18: Queue length StdDev for different Qeq thresholds in switch 3 with the multi-link topology

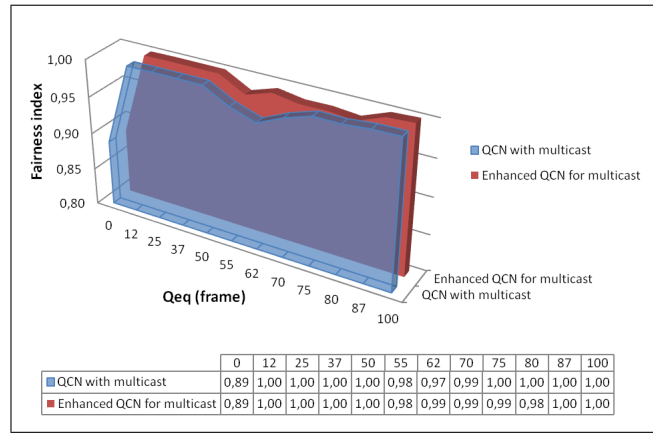


Figure 19: Fairness index of G1 with the multi-link topology

We note that the fairness indices are similar between the two QCN schemes for the majority values of the Qeq threshold.

5 Discussion

The proposed solution to enhance the QCN for multicast traffic is able to reduce the feedback overhead notably compared to an initial solution of the QCN with multicast traffic.

The proposed solution to enhance the QCN for multicast traffic shows a performance degradation in terms of loss rate and stability when the Qeq threshold value increases. Indeed, there is already a low feedback rate with the QCN with multicast traffic when the Qeq threshold value increases. This is due to the fact that high Qeq threshold values are exceeded less frequently than low values of the Qeq threshold. In addition, when the feedback frames are generated to notify sources about the state of congestion, it is generally late because there is not enough safety margin before reaching to the queue capacity. Therefore we note that high Qeq threshold values are characterized by high frame loss. With the selective technique of the representative that the enhanced QCN for multicast traffic exploits, one major consequence can be

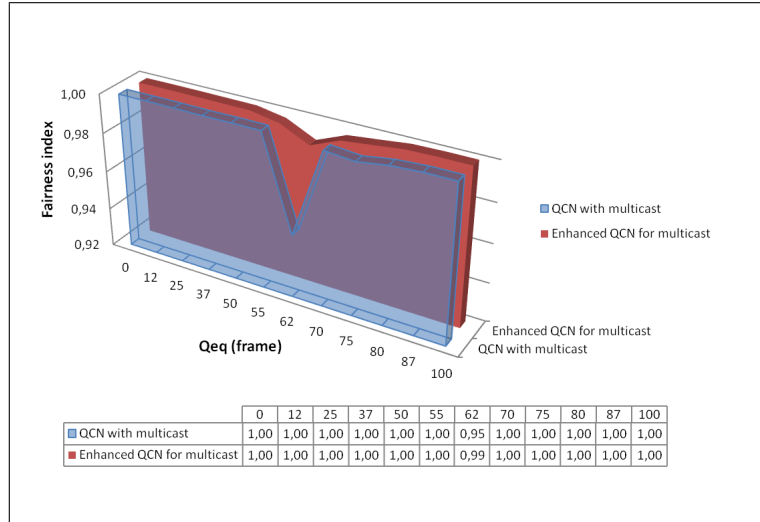


Figure 20: Fairness index of G2 with the multi-link topology

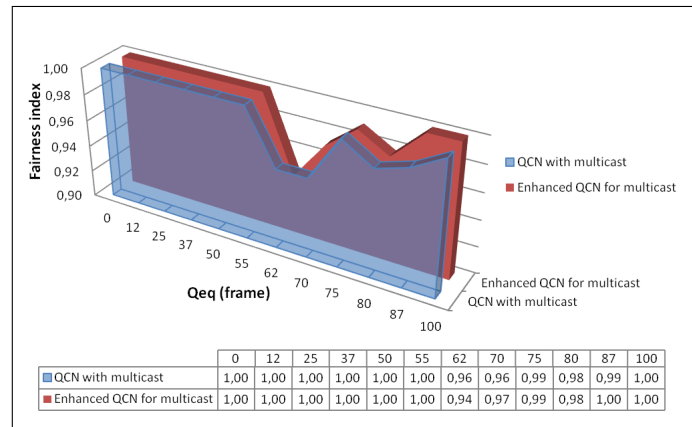


Figure 21: Fairness index of G3 with the multi-link topology

emphasized: the sources receive a poor notification information of congestion defined by a low number of feedback frames. Consequently, this does not allow the RP to adjust the sending rate adequately, and then both the loss rate and stability are not well maintained in the case of the enhanced QCN for multicast traffic when the Qeq threshold values are increased.

However, the standard recommendation of the QCN [9], [22] requires to chose a low Qeq value (25% of the queue capacity) to provide a safety margin for burst arrivals of new flows and to decrease queue delay. With such a recommendation, the proposed solution to enhance QCN for multicast traffic gives satisfactory results in terms of rate variation, feedback overhead, frame loss rate and scalability. However, stability and fairness performance results seem to be similar for both QCN schemes.

Table 1 shows a comparison summary between the QCN with multicast traffic and the enhanced QCN for multicast traffic.

6 Conclusion

The Quantized Congestion Notification (QCN) is a Layer 2 congestion control scheme for Carrier Ethernet data center network. Its purpose is to prevent the queue from building up to the point at which a frame has to be dropped causing it to use a Qeq threshold parameter. The QCN adjusts the source sending rate traffic according to the received feedback frames generated by the congestion point (i.e., the switch). Upon detecting congestion, the source needs to perform the appropriate transmission rate adjustments. When the source receives a feedback frame, it decreases its sending rate. Then, it undergoes successive rate increase phases: Fast Recovery (FR), Active Increase (AI) and Hyper Active Increase (HAI).

Due to the dearth of studies for the QCN in the case of multicast traffic, our aim in this paper was to study the QCN in the case of multicast traffic, and to propose an enhancement for the QCN to better handle multicast traffic. Our proposed schema for multicast traffic can contribute significantly to limiting the number of feedback frames without adding complex operations to the standard QCN.

The enhanced QCN for multicast traffic is inspired from the representative technique. This schema proposes to select a CP as a congestion representative among all the existing potential CPs in the path of the multicast traffic. The choice of a CP as a representative is updated each time a potential CP computes the greatest congestion measure $|Fb|$. The proposition does not require a receiver to source RTT computation to address the feedback implosion problem. Neither does it require knowledge of group membership or network topology. This could fit the standard congestion control scheme QCN as it does not involve additional information that the QCN standard does not require.

This paper evaluated, through simulations, the QCN performance for multicast traffic in terms of rate variation, feedback overhead, loss rate, stability, fairness and scalability. This paper also compares between the QCN with multicast traffic and the enhanced QCN for multicast traffic. We carried out traffic simulations for different Qeq threshold values. It appears from our findings that the enhanced QCN for multicast traffic has better performance than the QCN with multicast traffic in terms of rate

variation, feedback overhead, loss rate and scalability.

Our future work should focus on the study of the QCN for multicast traffic in heterogeneous network links parameter scenarios. This should be useful in the ongoing efforts to expand the deployment of the Carrier Ethernet network.

References

- [1] R. Fu, Y. Wang, and M.S. Berger. Carrier ethernet network control plane based on the next generation network. *IEEE Kaleidoscope Academic Conference Innovations in NGN: Future Network and Services*, pages 198–293, May 2008.
- [2] R. Fu, M.S. Berger, Y. Zheng, L. Brewka, and H. Wessing. Next generation network based carrier ethernet test bed for IPTV traffic. *IEEE EUROCON*, pages 1781–1787, May 2009.
- [3] A. Greenberg, P. Lahiri, D. A. Maltz, P. Patel, and S. Sengupta. Towards a next generation data center architecture: scalability and commoditization. *ACM workshop on Programmable Routers for Extensible Services of Tomorrow*, pages 57–62, 2008.
- [4] MEF Technical Specification. *EVC Ethernet Services Definitions Phase 3*. MEF 6.2, 2014.
- [5] MEF Technical Specification. *Carrier Ethernet Services for Cloud Implementation Agreement*. MEF 47, 2014.
- [6] D. Cai, A. Wielosz, and S. Wei. Evolve carrier ethernet architecture with SDN and segment routing. *IEEE 15th International Symposium on a World of Wireless, Mobile and Multimedia Networks (WoWMoM)*, pages 1–6, June 2014.
- [7] H. Mliki, L. Chaari, and L. Kamoun. A comprehensive survey on carrier ethernet congestion management mechanism. *Elsevier Journal of Network and Computer Applications*, 47:107–130, January 2015.
- [8] W. Jiang, F. Ren, Y. Wu, C. Lin, and I. Stojmenovic. Analysis of backward congestion notification with delay for enhanced ethernet networks. *IEEE Transactions on Computers*, 63:2674–2684, November 2014.
- [9] IEEE Computer Society. *IEEE Std 802.1 Qau Amendment 13: Congestion Notification. Local and Metropolitan area Networks-Virtual Bridged Local Area Networks*, 2010.
- [10] W. Jiang, F. Ren, and C. Lin. Phase plane analysis of quantized congestion notification for data center ethernet. *IEEE/ACM Transactions on Networking*, 23:1–14, February 2015.
- [11] M. Alizadeh, B. Atikoglu, A. Kabbani, A. Lakshmikantha, P. Rong, B. Prabhakar, and M. Seaman. Data center transport mechanisms: Congestion control theory and IEEE standardization. *Annual Allerton conference on Communication, Control and Computing*, pages 1270–1277, 2008.

- [12] Yang-Hua Chu, Sanjay G. Rao, and Hui Zhang. A case for end system multicast (keynote address). *ACM SIGMETRICS International conference on Measurement and Modeling of Computer Systems*, pages 1–12, 2000.
- [13] K. Sreerekha and V.K. Kiran. Mitigating incast congestion with LTTP for many to one communication in data centers. *IEEE International Conference on Innovations in Information, Embedded and Communication Systems (ICIIECS)*, pages 1–6, March 2015.
- [14] S.A. Pistirica, O. Poncea, and M.C Caraman. QCN based dynamically load balancing: QCN weighted flow queue ranking. *IEEE International Conference on Control Systems and Computer Science (CSCS)*, pages 197–204, May 2015.
- [15] A. Kabbani, M. Alizadeh, M. Yasuda and P. Rong, and B. Prabhakar. AF-QCN: Approximate fairness with quantized congestion notification for multi-tenanted data centers. *IEEE Annual Symposium on High Performance Interconnects (HOTI)*, pages 58–65, August 2010.
- [16] D. Crisan, A.S. Anghel, R. Birke, C. Minkenberg, and M. Gusat. Short and fat: TCP performance in CEE datacenter networks. *IEEE Annual Symposium on High Performance Interconnects (HOTI)*, pages 43–50, August 2011.
- [17] Y. Zhang and N. Ansari. On mitigating TCP incast in data center networks. *Proceedings IEEE INFOCOM*, pages 51–55, April 2011.
- [18] Y. Hayashi, H. Itsumi, and M. Yamamoto. Improving fairness of quantized congestion notification for data center ethernet networks. *IEEE International Conference on Distributed Computing Systems Workshops (ICDCSW)*, pages 20–25, June 2011.
- [19] Z. Yibo, E. Hagai, F. Daniel, G. Chuanxiong, L. Marina, L. Yehonatan, P. Jitendra, R. Shachar, Y. M. Haj, and Z. Ming. Congestion control for large-scale RDMA deployments. *SIGCOMM Proceedings ACM Conference on Special Interest Group on Data Communication*, pages 523–536, October 2015.
- [20] R. Mittal, T. Lam, N. Dukkupati, E. Blem, H. Wassel, M. Ghobadi, A. Vahdat, Y. Wang, D. Wetherall, and D. Zats. Timely: Rtt-based congestion control for the datacenter. *SIGCOMM Proceedings of the ACM Conference on Special Interest Group on Data Communication*, pages 537–550, October 2015.
- [21] H. Mliki, L. Chaari, L. Kamoun, and B. Cousin. Performance evaluation of legacy QCN for multicast and multiple unicast traffic transmission. *International Journal of Network Management*, pages 1–25, March 2016.
- [22] A. S Anghel, C. Basso, R. Birke, D. Crisan, M. Gusat, K. G Kamble, and C. J. Minkenberg. Quantized congestion notification (QCN) extension to explicit congestion notification for transport-based end-to-end congestion notification. <http://www.freepatentsonline.com/y2015/0188820.html>, July 2015.
- [23] D. DeLucia and K. Obraczka. Multicast feedback suppression using representatives. in: *proceeding of IEEE INFOCOM .Sixteenth Annual Joint Conference of the IEEE*

Computer and Communications Societies. Driving the Information Revolution, 2:463–470, April 1997.

- [24] L. Rizzo. Pgmcc: a TCP-friendly single-rate multicast congestion control scheme. *ACM SIGCOMM Computer Communication Review*, 30:17–28, October 2000.
- [25] J. Widmer and M. Handley. TCP-friendly multicast congestion control (tfmcc): Protocol specification. *RFC 4654*, August 2006.
- [26] J.P. Macker and R.B. Adamson. A TCP friendly, rate-based mechanism for NACK-oriented reliable multicast congestion control. *IEEE GLOBECOM*, 3:1620–1625, 2001.
- [27] J. Padhye, V. Firoiu, D. Towsley, and J. Kurose. Modeling TCP throughput: a simple model and its empirical validation. *in: Proceeding of ACM SIGCOMM conference on Applications, technologies, architectures, and protocols for computer communication*, pages 303–314, 1998.
- [28] R. Jain, W. Hawe, and D. Chiu. A quantitative measure of fairness and discrimination for resource allocation in shared computer systems. *DEC-TR-301*, pages 1–38, September 1984.
- [29] T. Bonald, L. Massoulié, A. Proutiere, and J. Virtamo. A queueing analysis of max-min fairness, proportional fairness and balanced fairness. *Queueing Systems*, 53:65–84, June 2006.

Table 1: *A comparison between the QCN with multicast traffic and the enhanced QCN for multicast traffic schemes*

Performance criteria	QCN with multicast scheme	Enhanced QCN for multicast scheme
The RP rate variation	- Less traffic rate than the enhanced QCN for multicast traffic scheme.	- Improves the average traffic rate at the RP.
Feedback overhead	- The feedback rate is higher than the enhanced QCN for multicast traffic scheme. - The more the Qeq threshold decreases, the more the queue generates feedback frames.	- The feedback rate is reduced.
Loss rate	- The loss rate is higher than the enhanced QCN for multicast traffic scheme. - As the Qeq threshold increases, the loss rate also increases.	- The loss rate is reduced when the Qeq threshold values are low.
Stability	- Better stability when the Qeq threshold values are increased.	- Better stability when the Qeq threshold values are decreased.
Fairness	- Similar fairness.	
Scalability	- The performance criteria degrade in terms of feedback overhead and loss rate.	- Decreases the feedback overhead and the loss rate.